

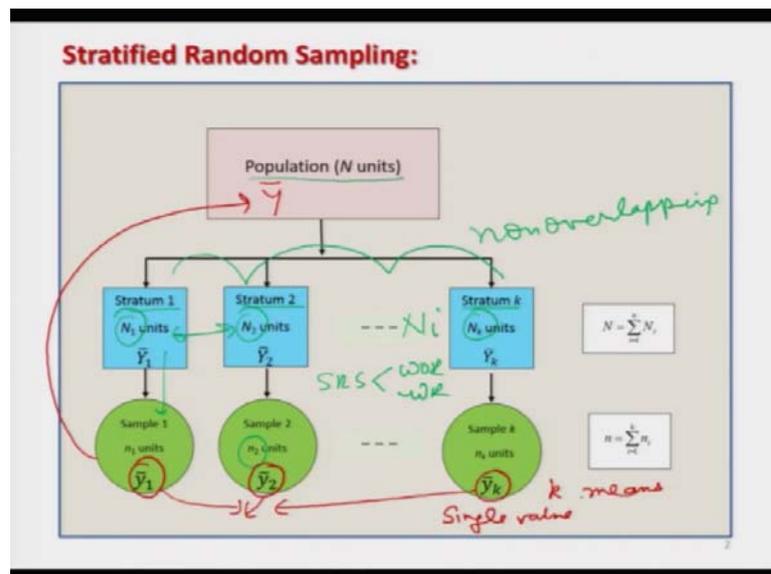
Essentials of Data Science with R Software - 2
Sampling Theory and Linear Regression Analysis
Prof. Shalabh
Department of Mathematics and Statistics
Indian Institute of Technology, Kanpur

Sampling Theory with R Software
Lecture - 29
Stratified Random Sampling
Estimation of Population Mean, Population Variance and Confidence Interval

Hello friends, welcome to the course Essentials of Data Science with R Software-2, where we are trying to learn the topics of Sampling Theory and Linear Regression Analysis. In this module on Sampling Theory with R Software, we are going to continue with our chapter Stratified Random Sampling.

So, in the last lecture, we started a discussion on a stratified random sampling, and I had given you a fair idea that under what type of conditions the stratified random sampling can be used. Now, in this lecture I will move forward, and we will see how we can estimate the population mean and population variance, and how they can be used in the construction of confidence interval ok. So, let us begin our lecture.

(Refer Slide Time: 01:12)



But let me try to take the slide of the last lecture to explain you some more things. So, you may recall that we had used these symbols that we have a population of size capital

N which is divided into say k different non-overlapping they are non-overlapping groups or they are called as strata or this is strata number 1, strata number 2, and stratum k .

So, all those units of population are divided in this case strata. They have been divided in such a way such that within group, they are the sampling codes are homogeneous with respect to the characters under study. And when we try to say between groups between strata they are heterogeneous as far as possible. So, remember as so just remember within homogeneous as far as possible and between heterogeneous as far as possible right.

So, now, these capital N units are divided in capital N_1 units, N_2 units, ..., N_k units. So, capital N , N_i , this is going to denote the number of units in the i th strata and from there we can use the SRS either WOR or say WR to draw the samples. And from the stratum 1, we draw a sample of size small n_1 number of units from capital N_1 number of units. And in the second sample, we try to draw see small n_2 number of units from population of size capital N_2 number of units and so on.

Now, the issue comes here. Based on the sample number 1, we have obtained here the sample mean \bar{y}_1 ; based on sample number 2, we have computed the sample mean \bar{y}_2 ; and similarly from the k th sample, we have computed the k th sample mean. But what is your objective?

Your objective is not to consider this k different sample means, but you are more interested in getting a single value; that means, you want to combine all these $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_k$. And you want to combine them in such a way such that their combination gives you a good estimator of the population mean which exist somewhere here.

So, one option to combine the different sample mean is to consider the weighted arithmetic mean. Well, there can be other approaches also. But this idea is coming from the statistical inference that if you try to consider the values of different weighted arithmetic mean, and yeah means you try to choose those weights in such a way so that they are computable and they and the resultant estimator gives you a good value of the population parameter.

(Refer Slide Time: 04:44)

Stratified Random Sampling:

We use the following symbols and notations:

N : Population size

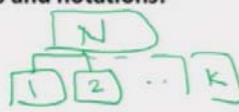
K : Number of strata

N_i : Number of sampling units in i^{th} strata

$N = \sum_{i=1}^k N_i$ Total population size

n_i : Numbers of sampling units to be drawn from i^{th} stratum.

$n = \sum_{i=1}^k n_i$: Total sample size



So, with this objective we move forward and first let us try to fix our notations and symbols. So, I am going to use here capital N to denote the population size. So, there are capital N number of sampling units. And from these capital N number of units, we will create capital K number of strata, that means, this population of size capital N is going to be divided into say difference strata 1, number strata 1, strata 2, up to here strata k.

And what is here capital N_i ? N_i is the number of sampling units in the i^{th} strata. And if you try to sum all the strata size sizes, then it will give us the total sample population size. And similarly from each of the strata, you try to draw the sample. So, this small n_i is the number of sampling units to be drawn from the i^{th} stratum. And the sum of all n_1, n_2, \dots, n_k , will give us the total sample size, so that is going to be our notation during the stratified sampling lecture.

(Refer Slide Time: 05:54)

Stratified Random Sampling:

Let

- Y : characteristic under study,
- y_{ij} : value of j^{th} unit in i^{th} stratum $j = 1, 2, \dots, n_i, i = 1, 2, \dots, k$,
- $\bar{Y}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} y_{ij}$: population mean of i^{th} stratum, $j = 1, 2, \dots, n_i, i = 1, 2, \dots, k$,
- $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$: sample mean of units from i^{th} stratum
- $\bar{Y} = \frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_i = \sum_{i=1}^k w_i \bar{Y}_i$: population mean where $w_i = \frac{N_i}{N}$.

And let this capital Y denote the characteristic under study, for example, say height or weight or age or anything else. And small y_{ij} is going to denote the value of the j th unit in the i th stratum. This means here what? For example, if you have here a population, you are divided you are dividing this into some samples.

So, now there is some unit say Y_1 some unit Y_2 . Now, this unit may be classified here, say for example, say in stratum number 1 and y_2 may be divided into say a stratum number 2. So, once it comes here from there, now this unit will have two addresses. First address is that when it is included in the stratum, whether this is the first unit to be to be included in the stratum, or second unit or say j th unit.

So, this j will go from 1 to n_i means n_i is the number of observations which are drawn from the i th strata right. So, here is our sample, somewhere here it will be y_{ij} . And then the second address will be whether this Y_1 is going to say stratum 1, stratum 2 or say stratum k , right.

So, y_{ij} is going to denote the value of the j th unit in the i th stratum. And now we try to find out the mean of all the units in the i th stratum. So, the mean of i th stratum that is the population mean of i th stratum will be denoted by \bar{Y}_i , and this is going to be the sum of all the units which are present in the i th stratum right.

So, this is denoted by $\frac{1}{N_i} \sum_{j=1}^{N_i} y_{ij}$. And similarly whatever are the units which are drawn

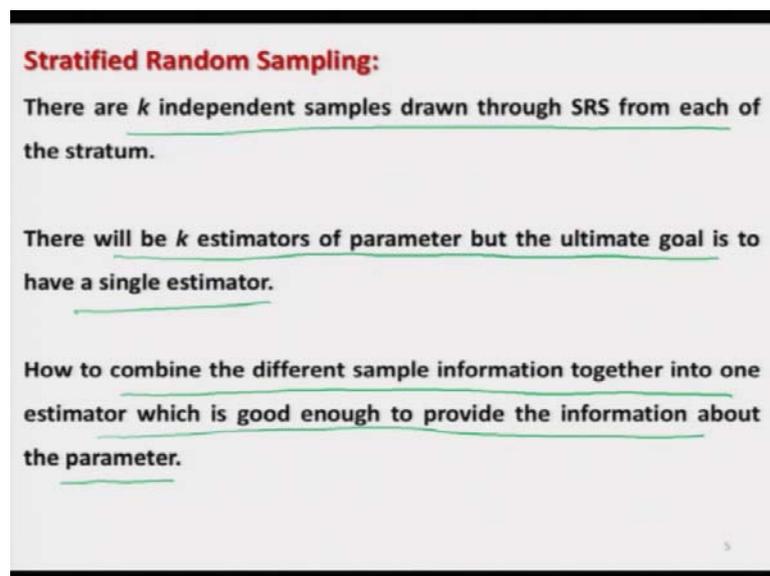
from the i th stratum, they are denoted by \bar{y}_i say $\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$. So, this is the sample mean of the units which are drawn from the i th stratum.

And then \bar{Y} , \bar{Y} is the population mean which is the mean of all the units in the population capital Y_1, Y_2, \dots, Y_n ; but here in this case we are not considering the simple arithmetic mean to indicate the \bar{Y} , but we are denoting here a weighted arithmetic mean.

So, this weighted arithmetic mean is given by $\frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_i$, right. So, I can write down this

quantity as here say $\frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_i$. So, this N_i/N , I am trying to denote by here w_i , right ok.

(Refer Slide Time: 09:36)



Stratified Random Sampling:

- There are k independent samples drawn through SRS from each of the stratum.
- There will be k estimators of parameter but the ultimate goal is to have a single estimator.
- How to combine the different sample information together into one estimator which is good enough to provide the information about the parameter.

So, now so from every strata, we try to draw a sample independently. So, ultimately, what we will have? We will have k independent sample which are drawn by simple random sampling from each of the stratum. So, and each independent sample will yield one value of the sample mean. So, there will be k estimator of the parameter population mean but our ultimate objective and goal is to have a single estimator.

So, the question comes here, how to combine those different k sample information together into one estimator which is good enough to provide the information about the parameter? Here at this stage, I would like to inform all of you that when we are trying to deal in statistics, there is a topic what we say this combining different sample information that is itself a topic in statistics, where we where people try to address that how different type of sample information can be combined together to give a good outcome.

(Refer Slide Time: 10:48)

Estimation of population mean:
 In case of stratified sampling, the population mean is defined as the weighted arithmetic mean of stratum means where the weights are provided in terms of strata sizes.

$$\bar{Y} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$$

Find the sample mean of the units drawn from each stratum.
 Find their weighted mean, called as stratified mean.

Use stratified mean as to estimate the population mean as

$$\rightarrow \bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i \quad w_i = \frac{N_i}{N} \quad \left(w_i = \frac{n_i}{n} \right) ?$$

So, in the case of stratified sampling as we have discussed the population mean is defined as the weighted arithmetic mean of stratum means right. So, now, these weights are actually known to us. So, this capital N_i is known to us, capital N is known to us. So, the weights w_i 's are known to us.

Now, we try to find out the sample mean of all the units which are drawn from each stratum, and we try to find out their weighted arithmetic mean. And this is called as stratified mean. This stratified mean is defined here like this $\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$, right.

So, here the weights are obtained using the size of the stratum N_i divided by N . Well, there is a there is always a question from the students at this stage why are we using the

size of the stratum, why not are we using the weights as say size of the sample say small n_i upon say n . So, I will try to address this thing in the next couple of slides, ok.

So, now we have defined that this is our estimator \bar{y}_{st} that is going to estimate the population mean. The story behind this thing that why cannot we use the weighted mean of sample median, sample mode or something else that reason goes to the say first couple of lectures when we started the discussion on simple random sampling right, so I will not repeat it here again.

(Refer Slide Time: 12:57)

Unbiased estimator of population mean

Since the sample in each stratum is drawn by SRS, so

$$E(\bar{y}_i) = \bar{Y}_i \quad N_i \quad n_i \quad E(\bar{y}_i) = \bar{Y}_i$$

$$\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$$

$$E(\bar{y}_{st}) = \frac{1}{N} \sum_{i=1}^k N_i E(\bar{y}_i) \rightarrow \bar{Y}_i$$

$$= \frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_i$$

$$= \bar{Y}$$

\bar{y}_{st} is an unbiased estimator of \bar{Y}

But now we have got here an estimator which we are proposing that which is to be used to find out the population mean, right. So, we would like to study its property whether this estimator is good or bad; and based on that we will make our recommendation. Now, before we try to understand or before we try to do the algebra in this topic, I would like to just tell you something, you see we already have done all the algebra very carefully during the simple random sampling.

I have shown you how to prove unbiasedness, how to prove how to find out the variances or their estimates. And since here we are using simple random sampling, so I can use those results here directly. For example, when I say I have a stratum of size capital N_i from there I am trying to draw here a small sample of size small n_i by SRS, it can be WR or WOR.

That means, if I ask you that what will happen if I try to consider this i th sample mean as my sample mean, and this i th stratum as my population? Obviously we are trying to draw the sample by simple random sampling. So, we already have proved that this is going to be an unbiased estimator. So, I do not need to prove it again. The only thing I have to prove that if I try to combine them using the weights w_i , then what will happen?

Similarly, when we try to find out the variance, we already have proved that if you try to take this thing as a sample which is drawn from the i th stratum, then, and if you try to consider this stratum as the population then we know what is the expression and how we can find out the expression for the variance of \bar{y} that is the sample mean when we are trying to draw a sample size sample of size small n_i from a population of size capital N_i .

So, I do not need to prove them again, and that was the reason that I had given you details in those cases. But it does not mean that I am going to give you here all the proofs. As I declared in the beginning itself, I will try to give you those proofs which I can handle here easily, and which are essential for the development of the course and to explain you the concept.

There are many, many results for example. In this stratified sampling also which are important. So, I will tell you those important result and some results I will refer you to some other books or courses. So, we will try to do in this way right. So, let us now start our discussion. Since the sample in each of the sample is being drawn by simple random sampling, so obviously if you try, so if you have learned that if you have a population of size capital N and if you draw the sample of size small n , then expected value of \bar{y} here is capital \bar{Y} .

Now, what I am asking you now your population size is N_i – capital N_i , and your sample size is small n_i , based on that your sample mean is \bar{y}_i . So, do not you think that this is automatically going to be an unbiased estimator of capital \bar{Y}_i ? Answer is yes, so why not to use this information directly and intelligently.

So, now I try to find out the expected value of \bar{y}_{st} . So, that is going to be $\frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$, and this expectations operator which comes inside the summation sign. So, now, I know this expected value of \bar{y}_i is nothing but your \bar{Y}_i . So, I try to substitute it here.

And if you try to see this quantity $\frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$, this is nothing but your \bar{Y} which is an unbiased estimator of \bar{Y} . So, in this case, \bar{y}_{st} is an unbiased estimator of \bar{Y} . So, this \bar{y}_{st} has got nice property that it is unbiased. Now, what is the meaning of unbiased? So, now we have got good estimator.

(Refer Slide Time: 17:44)

Biased estimator of population mean

Since the sample in each stratum is drawn by SRS, so

$$E(\bar{y}_i) = \bar{Y}_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^k n_i \bar{y}_i$$

$$E(\bar{y}) = \frac{1}{n} \sum_{i=1}^k n_i E(\bar{y}_i)$$

$$= \frac{1}{n} \sum_{i=1}^k n_i \bar{Y}_i$$

$$\neq \frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_i$$

$$\neq \bar{Y} \rightarrow \bar{y} \text{ is a biased estimator of } \bar{Y}$$

$$\bar{y}_{st} = \sum_{i=1}^k w_i \bar{y}_i$$

$$w_i = \frac{N_i}{N}$$

$$w_i^* = \frac{n_i}{n}$$

$$\bar{y} = \sum_{i=1}^k w_i^* \bar{y}_i$$

$$= \sum_{i=1}^k \frac{n_i}{n} \bar{y}_i$$

But now first I try to address the question which I raised earlier that in this case, we have obtained an estimator of population mean by considering the weighted arithmetic mean \bar{y}_i where your weights are given by capital N_i upon N using the stratum size. Now, suppose I say that instead of stratum size, I try to use here sample size. Say and I write down the new weights as a small n_i upon n .

So, in this case, if I try to define here a \bar{y} , this will become here $\sum_{i=1}^k w_i \bar{Y}_i$. So, this

becomes here $\frac{1}{n} \sum_{i=1}^k n_i \bar{Y}_i$. So, what will now happen to this estimator that we have to see.

So, this is the same estimator which I have written here ok.

So, now, once I try to check the unbiased behavior of this estimator, we know that expected value of \bar{y}_i will be equal to capital \bar{Y}_i . Now, once I take this expectation sign operator inside the summation sign while finding out the expected value of \bar{y} , this expected value of \bar{y}_i will be then nothing but your capital \bar{Y}_i .

And then you can see here this expression is not the same as the value of \bar{Y} which you have defined. Now, in case, if you argue that I can define the population mean by this thing $\frac{1}{n} \sum_{i=1}^k n_i \bar{Y}_i$ possibly this is not correct, because you are taking only here n_i

observation which are not representing the population, the population has capital N_i observation in every strata. And you are choosing only the selected number of observation from average strata to indicate the population mean which is itself contradictory.

So, expected value of \bar{y} does not comes out to be \bar{Y} . So, I can say here \bar{y}_i is a biased estimator of population mean \bar{Y} , and that is the reason we do not consider \bar{y} , but we consider here \bar{y}_{st} as an as a good estimator for estimating the population mean.

(Refer Slide Time: 20:44)

Variance of stratum mean:

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 Var(\bar{y}_i) + \sum_{i \neq j=1}^k \sum_{j=1}^{n_i} w_i w_j Cov(\bar{y}_i, \bar{y}_j)$$

$\bar{y}_{st} = \sum_{i=1}^k \omega_i \bar{y}_i$

Since all the samples have been drawn independently from each of the strata by SRSWOR, so

$$Cov(\bar{y}_i, \bar{y}_j) = 0, \quad i \neq j$$

$$Var(\bar{y}_i) = \frac{N_i - n_i}{N_i n_i} S_i^2$$

where

$$S_i^2 = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (Y_{ij} - \bar{y}_i)^2$$

Handwritten notes: N_i (circled), n_i (circled), SRSWOR, $\frac{N_i - n_i}{N_i n_i} S_i^2$, $S_i^2 = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (Y_{ij} - \bar{y}_i)^2$

Now, I try to find out the variance of this \bar{y}_{st} . So, remember one thing, I will just write

down here for your convenience $\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$. And if you try to take out try to operate

the variance operator and if you try to find out variance of \bar{y}_{st} , this will come out to

$$be Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 Var(\bar{y}_i) + \sum_{i \neq j=1}^k \sum_{j=1}^{n_i} w_i w_j Cov(\bar{y}_i, \bar{y}_j).$$

So, now you can see here this is the covariance between the ith and jth sample and ith and jth sample have been drawn independently. In fact, all the sample have been drawn independently of each other from all the strata. So, this covariance will become 0. So, that is what I am trying to write down here that covariance between \bar{y}_i and \bar{y}_j will be 0.

And when you are trying to draw a sample of size n_i from a population of size capital N_i by suppose SRS WOR, then we know that earlier the variance was given by $\frac{N - n}{Nn} S^2$.

Now, I can make it here because this is coming from ith strata this will become N_i this will become small n_i , this is here capital N_i small n_i and S^2 will become S_i^2 , right.

So, this is exactly what I am trying to do here that is the variance of \bar{y}_i from the i th strata that the variance of the sample mean from the i th strata will be denoted by this quantity right. And what is here S^2 or S_i^2 ? You if you remember S square was earlier the variance of the population which is in the capital N_i .

So, now I have to rewrite this S square as earlier it was $S^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2$. Now, I will try to rewrite this definition for the i th strata. So, this S^2 becomes S_i^2 and capital N becomes n_i and this Y_i becomes Y_{ij} and this \bar{Y} becomes \bar{Y}_i . So, this is here the definition of S_i^2 .

(Refer Slide Time: 23:37)

Variance of stratum mean:

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 Var(\bar{y}_i) + \sum_{i \neq j=1}^k \sum_{j=1}^{n_i} w_i w_j Cov(\bar{y}_i, \bar{y}_j)$$

Thus

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \frac{N_i - n_i}{N_i n_i} S_i^2$$

→ SRSWOR ←

$$= \sum_{i=1}^k w_i^2 \left(1 - \frac{n_i}{N_i} \right) \frac{S_i^2}{n_i}$$

$$= \sum_{i=1}^k w_i^2 \left(1 - \frac{n_i}{N_i} \right) \frac{S_i^2}{n_i}$$

So, now you can see here you have obtained all these expression. And you simply have to use this expression over here, and you will get the expression right. So, I can just substitute that expression over there and you can see here just substituting it here, I get the variance of \bar{y}_{st} under SRSWOR like this right. One thing you can observe here that this same quantity I can write down here as say w_i square $1 - \frac{n_i}{N_i}$ upon capital N_i , if you remember we had denoted a small n by capital N as a sampling fraction. So, yes, I can write in this particular way also, right ok.

(Refer Slide Time: 24:20)

How to construct strata:

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \left(1 - \frac{n_i}{N}\right) \frac{S_i^2}{n_i}$$

$w_i = \frac{N_i}{N}$

Handwritten notes:
 N_i is fixed
 S_i^2 is left
 Small
 S_i^2 is small

Variance is small when S_i^2 is small.

This suggests how to construct the strata.

If S_i^2 is small for all $i = 1, 2, \dots, k$, then variance will also be small.

That is why the strata are constructed such that they are within homogeneous, i.e., S_i^2 is small and among heterogeneous.

Now, at this address, I will try to show you that why it is recommended to construct the strata in such a way such that they are within homogeneous and between heterogeneous. You can see here in this expression, this is the variance of \bar{y}_{st} under SRSWOR. This quantity here is fixed that you cannot control because w_i is N_i upon N , right.

Now, there are only two things remaining which you have not controlled capital N_i – it is not in your control that you means once you have created the groups the stratum size is fixed. Now, there are two things here n_i and S_i^2 . At this moment, for a while we are assuming that n_i is fixed somehow means you have followed a rule and you have obtained the sample size.

So, the only quantity which is left to us is S_i^2 . . . Although I will also try to address that this small n_i can also be controlled, but that will come later. At this moment, I am assuming that this is fixed those who have knowledge of stratified sampling, they can imagine that the sample size can be obtained by different types of allocations equal allocation, proportional allocation, Neyman allocation and there are various things, but at this moment I am assuming that they are fixed.

So, now only S_i^2 is left in this quantity which is to be controlled. Obviously, now if this S_i^2 which is small then it will make the variance of \bar{y}_{st} also smaller. So, I have written

here very clearly if S_i^2 is small for all i 's means the capital S_i^2 for first strata, second strata, k th strata if all of them are say small, then the variance of \bar{y}_{st} will also be small

Note that if one of the S_i^2 is very very large possibly that will also disturb. So, ideally I would say that if all S_i^2 s are smaller, then variance of \bar{y}_{st} will also be smaller. Now, what does this mean? When I say S_i^2 is small, what does this mean? Try to think try to take a pause of this video, and try to think for a couple of minutes and then start the video right.

So, when I say S_i^2 is small, that means, the variability within the i th group is smaller, then the variance of \bar{y}_{st} will be small. So, this is the reason that the strata are constructed in such a way such that the within group variability is as small as possible and then obviously the variability among the strata that will automatically become larger so that is the reason.

So, now, you can see here I have established that why you have considered the form of

\bar{y}_{st} to be like $\frac{1}{N} \sum_{i=1}^k N_i \bar{y}_i$, because this estimator has been constructed in such a way such

that it will give you a value of the population mean with the smaller variability. Yes, one can think about different other estimators also and that is one of the area of research in sampling theory where people try to find out better estimator who have got a smaller sampling variability, right, ok.

(Refer Slide Time: 29:01)

Unbiased estimation of variance of \bar{y}_{st} under SRSWOR:

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \left(\frac{N_i - n_i}{N_i} \right) \frac{S_i^2}{n_i}$$

Since the samples have been drawn by SRSWOR, so

$$E(s_i^2) = S_i^2 \rightarrow \text{WOR}$$

where $s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$

and $\widehat{Var}(\bar{y}_i) = \frac{N_i - n_i}{N_i n_i} s_i^2$

so $\widehat{Var}(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \widehat{Var}(\bar{y}_i)$

$$= \sum_{i=1}^k w_i^2 \left(\frac{N_i - n_i}{N_i n_i} \right) s_i^2$$

Find standard error by $+\sqrt{\widehat{Var}(\bar{y}_{st})}$.

Now, once you have obtained this variance of \bar{y}_{st} here, obviously, this depends on the population value capital S_i^2 . So, you cannot estimate it on the basis of given sample of data. So, now, we need to find out the estimate of variance of \bar{y}_{st} . But now it is not difficult, because you already have drawn the sample from the i th strata by SRSWOR.

And for that you already have proved during the lectures on simple random sampling that expected value of small s_i^2 capital S_i^2 in case of SRSWOR, where this S_i^2 is simply the value of s^2 in the i th sample. Once that is there, so I can replace this quantity capital S_i^2 by s_i^2 .

So, the estimate of variance of \bar{y}_i will come out to be directly from the SRSWOR that is

$$\frac{N_i - n_i}{N_i n_i} s_i^2$$

And I try to replace it here. So, the $\widehat{Var}(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \left(\frac{N_i - n_i}{N_i n_i} \right) s_i^2$. And this I

can replace here.

So, this quantity here $\sum_{i=1}^k w_i^2 \left(\frac{N_i - n_i}{N_i n_i} \right) s_i^2$. This is an estimate of the variance. So, once

you have a sample, you can also estimate the population variance of \bar{y}_{st} . And if you take the positive square root of this variance estimator, then you will get the standard error,

the issues which are related in defining the standard error that I already have discussed in the earlier lecture, so I will not repeat it again.

(Refer Slide Time: 31:05)

Unbiased estimation of variance of \bar{y}_{st} under SRSWR:

If the samples have been drawn by SRSWR, so

$$\bar{y}_{st} = \sum_{i=1}^k w_i \bar{y}_i \quad E(\bar{y}_{st}) = \sum_{i=1}^k w_i E(\bar{y}_i) = \bar{Y}$$

$$E(\bar{y}_{st}) = \bar{Y} \rightarrow \bar{y}_{st} \text{ is an unbiased estimator of } \bar{Y}$$

$$Var(\bar{y}_{st}) = \sum_{i=1}^k w_i^2 \left(\frac{N_i - 1}{N_i n_i} \right) S_i^2 = \sum_{i=1}^k w_i^2 \frac{\sigma_i^2}{n_i}$$

where $\sigma_i^2 = \frac{1}{N_i} \sum_{j=1}^{N_i} (y_{ij} - \bar{y}_i)^2$, $s_i^2 = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$

$$\widehat{Var}(\bar{y}_{st}) = \sum_{i=1}^k \frac{w_i^2 s_i^2}{n_i}$$

Find standard error by $+\sqrt{\widehat{Var}(\bar{y}_{st})}$

Handwritten notes on the right side of the slide:

- $\frac{N_i - 1}{N_i n_i} S_i^2$
- $(N_i - 1) S_i^2 = N_i \sigma_i^2$
- $E(S_i^2) = \sigma_i^2$

Now, I come to another aspect. Suppose, somebody decides that instead of for drawing the sample by SRSWOR from the i th stratum, suppose the decision is to draw the sample by SRSWR – Simple Random Sampling With Replacement. In that case, what will happen? So, let us try to find out the things.

So, in this case also the \bar{y}_{st} will remain the same. This is the same definition. But if you try to find out the expected value of here \bar{y}_{st} , this will become here $\sum_{i=1}^k w_i E(\bar{y}_i)$ and expected value of \bar{y}_i here is capital \bar{Y}_i . So, this quantity is nothing but your \bar{Y} .

So, under this case also \bar{y}_{st} remains an unbiased estimator of population mean capital \bar{Y} . So, there is no change in the unbiasedness character of \bar{y}_{st} . For finding of the variance of \bar{y}_{st} under WR in the case of SRSWR, we know that

$$\text{the } Var(\bar{y}) = \frac{N-n}{Nn} S^2.$$

So, once I try to use it here for the i th strata, that means, is a sample of size small n_i is obtained from a of the stratum of size capital N_i . So, this capital N will become N_i , and a small n will become small n_i , and S^2 will become the variance of the i th stratum.

So, this is what I have written here. And I simply replace the variance of \bar{y}_i by here this quantity. And you can see that this quantity is nothing but your $\sum_{i=1}^k w_i^2 \frac{\sigma_i^2}{n_i}$. Yeah, if you remember we had used this relation that $(N - 1)S^2 = N \sigma^2$, right.

And if you want to find out an n by estimator of variance of \bar{y}_{st} , so you may recall that we already had proved that expected value of s_i^2 or expected value of here s^2 is equal to σ^2 . So, in the case of stratified sampling, expected value of s_i^2 will become say σ_i^2 . So, I can simply replace here σ_i^2 by this s_i^2 .

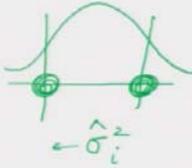
And you can see here an estimate of the variance of \bar{y}_{st} can be obtained by $\sum_{i=1}^k \frac{w_i^2 s_i^2}{n_i}$. And σ_i^2 can be defined exactly on the same line as we have defined capital S_i^2 . And if you want to find out the standard error, simply take the positive square root of this estimate of variance.

(Refer Slide Time: 34:01)

Confidence intervals of population mean:
 Assume \bar{y}_{st} is normally distributed, and $\sqrt{\widehat{Var}(\bar{y}_{st})}$ is well determined so that t can be read from normal distribution tables.

If only few degrees of freedom are provided by each stratum, then t values are obtained from the table of student's t -distribution. The confidence limits of \bar{Y} can be obtained as

$$\bar{y}_{st} \pm t \sqrt{\widehat{Var}(\bar{y}_{st})}$$



14

Now, I try to address the issue of confidence interval estimation of population mean in case of a stratified random sampling. Well, it is not so straight forward as in the case of simple random sampling. Because when we are trying to find out the confidence interval if you remember we need to find out here a critical value which are here on the x-axis corresponding to the distribution.

And usually the values of σ^2 or in case of stratified sampling the value of σ_i^2 are not known to us. So, we try to estimate them. And when we are trying to estimate them, then the properties of the \bar{y}_{st} changes and this change is quite different what happens in the case of simple random sampling, because we are trying to combine it together.

So, first of all we have to make an assumption that \bar{y}_{st} is normally distributed. And standard error of \bar{y}_{st} is well determined, so that the critical value say t_i am not really saying that this is the t distribution, but that I will show you that this is a critical value. This critical value can be read from the normal distribution table. And second thing is this if you remember the distribution of t also depends on the degrees of freedom.

So, in case if only a few degrees of freedom are provided by each stratum then t values can be obtained from the table of student's t-distribution. But this itself has an issue which I will try to address in the next slide. But suppose we assume that all these things are possible and under that situation an approximate confidence interval can be computed by this expression that $\bar{y}_{st} \pm t\sqrt{\widehat{Var}(\bar{y}_{st})}$.

refer Slide Time: 36:13)

Confidence intervals of population mean:

The distribution of $\sqrt{\text{Var}(\bar{y}_{st})}$ is generally complex.

Assume y_{ij} 's are normally distributed.

An approximate method of assigning an effective number of degrees of freedom n_e to $\sqrt{\text{Var}(\bar{y}_{st})}$ is

$$n_e = \frac{\left(\sum_{i=1}^k g_i s_i^2 \right)^2}{\sum_{i=1}^k \frac{g_i^2 s_i^4}{n_i - 1}}$$

$$g_i = \frac{N_i(N_i - n_i)}{n_i}$$

$$\text{Min}(n_i - 1) \leq n_e \leq \sum_{i=1}^k (n_i - 1)$$

But if you try to see this is not so straight forward. Why? Because the distribution of standard error of \bar{y}_{st} is not so simple. Well, that is the job of the statistician to find out all these things, but I am not giving you here the all the details. But as a data scientist you must know that what is the problem what is something, something is happening that you must know, so that you know that how much amount of mistakes you are going to make in your computation, and what are the violation because of which this error is coming, and how this violation can affect your final conclusion.

In case if you assume that all the observation y_{ij} 's they are coming from a normal distribution, then an approximate method for assigning the effective number of degrees of freedom to the distribution of standard error of \bar{y}_{st} , that means, why I am calling is effective number? Because you see the degrees of freedom are always some integer value.

And here the distribution of a standard error becomes too complex. So, it is not possible to get the exact number of degrees of freedom. So, we try to approximate it. And when we try to approximate it, the approximated value may not come out to be always an integer. So, that is why the value which we are going to estimate that is called an called as effective number of degrees of freedom.

So, the number of effective degrees of freedom which is denoted by a small n_e this is obtained by this expression, where this S_i^2 you know which is the variance from the i th sample and this quantity here g_i is defined here like this. And so if you try to compute here n_e , and then try to choose the value of n_e which is $\text{Min}(n_i - 1) \leq n_e \leq \sum_{i=1}^k (n_i - 1)$.

So, whatever you get here this value of n_e , try to choose this one and this value will lie between these two limits right ok. So, in this lecture, we have understood that why we have chosen this \bar{y}_{st} to estimate the population mean and how to obtain the standard errors which are needed in a real life data, and then after that how to construct the confidence interval.

One important aspect which you have learnt that whatever are the decisions, for example, the question which I raised in the beginning that whether the stratum or the strata are constructed in such a way that they are within homogeneous or between heterogeneous, these types of questions can be answered using the statistics, using the statistical concepts. And these are the rules which govern us that how we can construct an estimator for the population mean under different types of sampling scheme.

So, if you go for cluster sampling, if you go for systematic sampling, the estimated to estimate the population mean will be different they are not going to be always the same as say small y bar as in case of simple random sampling or \bar{y}_{st} in case of a stratified random sampling.

Well, these estimators are not coming from sky. Well, as a statistician this is the job that people try to construct different types of estimator, they try to find out their variances, and then they try to see which one is working better. And after making so much of efforts they come up finally with the good solution which I have given you here. So, now while constructing the strata there were two variables the sample size from the i th strata and S_i^2 .

So, I already have addressed that how to control the variability of \bar{y}_{st} using S_i^2 . But I assume that n_i is somehow known to us. In the next lecture, I will try to work with a small n_i . And we will try to see there are different ways which can give rise to different

types of sample sizes to be chosen from the strata. And then it is possible that the variance of \bar{y}_{st} will be further reduced. So, now you can see you have an estimator.

If you do not want to do anything just use it, but if you want to make the variability as small as possible, try to first construct the strata, and then after that you have one more option that you can decide that how much sample I can draw from every stratum, so that the final variability in my estimator is as small as possible. So, in the next turn, I will try to address the issue of sample selection.

So, you enjoy your lecture, you revise this, you study it, and I will take a leave till then, good bye.